

AD-A254 052



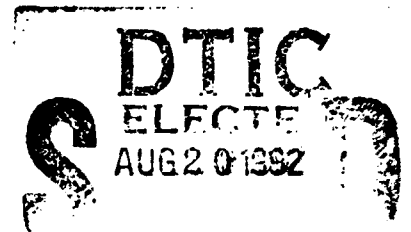
2

## Parity Declustering for Continuous Operation in Redundant Disk Arrays

Mark Holland† and Garth Gibson

April 1992

CMU-CS-92-130



School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

This document has been approved  
for public release and sale; its  
distribution is unlimited.

This material is based in part upon work supported by the National Science Foundation under grant number ECD-8907068, and in part upon work supported by the Defense Advanced Research Projects Agency (DoD) and monitored by DARPA/CMO under contract MDA972-90-C-0035. The government has certain rights in this material. Additional support was provided by an IBM Graduate Fellowship.

† Department of Electrical and Computer Engineering, Carnegie Mellon University.

92 8 19 86

92-23161



~~92 6 20 000~~

## Abstract

We describe and evaluate a strategy for declustering the parity encoding in a redundant disk array. This declustered parity organization balances cost against data reliability and performance during failure recovery in highly-available parity-based arrays for use in continuous-operation systems. It improves on standard parity organizations by reducing the additional load on surviving disks during the reconstruction of a failed disk's contents. This yields higher user throughput during recovery, and/or shorter recovery time.

We first demonstrate a software implementation of declustered parity based on balanced incomplete and complete block designs. This implementation is then evaluated using a disk array simulator under a highly concurrent workload comprised of small user accesses. We show that declustered parity penalizes user response time while a disk is being repaired (before and during its recovery) less than comparable non-declustered (RAID 5) organizations without any penalty to user response time in the fault-free state.

We then show that previously proposed modifications to a simple, single-sweep reconstruction algorithm further decrease user response times during recovery, but, contrary to previous suggestions, this may be achieved at the cost of slower recovery in many declustered parity arrays. This result arises from the simple model of disk access performance used in previous work, which did not consider throughput variations due to positioning delays.

DTIC QUALITY INSPECTED 3

Accession For	
NTIS CRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By <i>per lti.</i>	
Distribution/	
Availability	
Dist	Availability
<i>A-1</i>	

# 1. Introduction

Many applications, notably database and transaction processing, require both high throughput and high data availability from their storage subsystems. The most demanding of these applications require continuous operation, which in terms of a storage subsystem requires (1) the ability to satisfy all user requests for data even in the presence of a disk failure, and (2) the ability to reconstruct the contents of a failed disk onto a replacement disk, thereby restoring itself to a fault-free state. It is not enough to fulfill these two requirements with arbitrarily degraded performance; it is not unusual for a continuous operation application to suffer financial losses substantially larger than their investment in computing equipment if service is severely degraded for a prolonged period of time. Since the time necessary to reconstruct the contents of a failed disk is certainly minutes and possibly hours, we focus this paper on the performance of a continuous-operation storage subsystem during on-line failure recovery.

Redundant disk arrays, proposed for increasing input/output performance and for reducing the cost of high data reliability [Kim86, Livny87, Patterson88, Salem86], also offer an opportunity to achieve high data availability without sacrificing throughput goals. A single-failure-correcting redundant disk array consists of a set of disks, a mapping of user data to these disks that yields high throughput [Chen90b], and a mapping of a parity encoding for the array's data such that data lost when a disk fails can be recovered without taking the system off-line [Lee91].

Most single-failure-correcting disk arrays employ either *mirrored* or *parity-encoded* redundancy. In mirroring [Bitton88, Copeland89, Hsiao90], one or more duplicate copies of all data are stored on separate disks. In parity encoding [Kim86, Patterson88, Reddy89], popularized as Redundant Arrays of Inexpensive Disks (RAID), some subset of the physical blocks in the array are used to store a single-error-correction code (usually parity) computed over subsets of the data. Mirrored systems, while potentially able to deliver higher throughput than parity-based systems for some workloads [Chen90a, Gray90], increase cost by consuming much more disk capacity for redundancy. In this paper, we examine a parity-based scheme called *parity declustering*, which provides better performance during on-line failure recovery than more common RAID schemes, without the high capacity overhead of mirroring [Muntz90]<sup>1</sup>.

Our primary figures of merit in this paper are reconstruction time, which is the wallclock time taken to

reconstruct the contents of a failed disk after replacement, and user response time during reconstruction. Reconstruction time is important because it determines the length of time that the system operates at degraded performance, and because it is a significant contributor to the length of time that the system is vulnerable to data loss caused by a second failure. Given a fixed user throughput, contrasting user fault-free response time to the response time both before and during reconstruction gives us the measure of our system's performance degradation during failure recovery.

Section 2 of this paper describes our terminology and presents the declustered parity organization. Section 3 describes related studies, notably the introduction of declustering by Muntz and Lui [Muntz90], and explains the motivations behind our investigation. Section 4 presents our parity mapping, which was left as an open problem by Muntz and Lui. Section 5 gives a brief overview of our simulation environment and Sections 6 and 7 present brief analyses of the performance of a declustered array when it is fault-free, and when there is a failed disk but no replacement. Section 8 then covers reconstruction performance, contrasting single-thread and parallel reconstruction, and evaluating alternative reconstruction algorithms. Section 9 concludes the paper with a look at interesting topics for future work.

## 2. The Declustered Parity Layout Policy

Figure 2-1 illustrates the parity and data layout for a left-symmetric RAID 5 redundant disk array [Lee91]. A *data stripe unit*, or simply a *stripe unit* is defined as the minimum amount of contiguous user data allocated to one disk before any data is allocated to any other disk. A *parity stripe unit*, or simply a *parity unit*, is a block of parity information that is the size of a data stripe unit. The size of a stripe unit, called a *unit* for convenience, must be an integral number of sectors, and is often the minimum unit of update used by system software. A *parity stripe*<sup>2</sup> is the set of (data) stripe units over which a parity stripe

---

1. Muntz and Lui use the term *clustered* where we use the term *declustered*. Their use may be taken from "clustering" independent RAID's into a single array with the same parity overhead. Our use follows the earlier work of Copeland and Keller [Copeland89] where redundancy information is "declustered" over more than the minimal collection of disks.

2. A parity stripe associates stripe units that contribute to the same parity computation. This is not the same as a data stripe, often simply called a stripe, which might be defined as the maximum amount of contiguous user data that contains no more than one stripe unit from any disk. While most data and parity mappings are chosen to coincide for performance reasons, an array's parity stripe mapping is not dependent on its data stripe mapping.

Muntz and Lui [Muntz90] use the term *group* to denote what we call a parity stripe, but we avoid this usage as it conflicts with the Patterson, et. al. definition [Patterson88] as a set of disks, rather than a set of disk blocks.

unit is computed, plus the parity stripe unit itself. In Figure 2-1,  $D_{i,j}$  represents the  $j$ th stripe unit of parity stripe number  $i$ , and  $P_i$  represents the parity unit for parity stripe  $i$ . Parity units are distributed across the disks of the array to avoid the write bottleneck that would occur if a single disk contained all parity units.

Offset	DISK0	DISK1	DISK2	DISK3	DISK4
0	D0.0	D0.1	D0.2	D0.3	P0
1	D1.1	D1.2	D1.3	P1	D1.0
2	D2.2	D2.3	P2	D2.0	D2.1
3	D3.3	P3	D3.0	D3.1	D3.2
4	P4	D4.0	D4.1	D4.2	D4.3

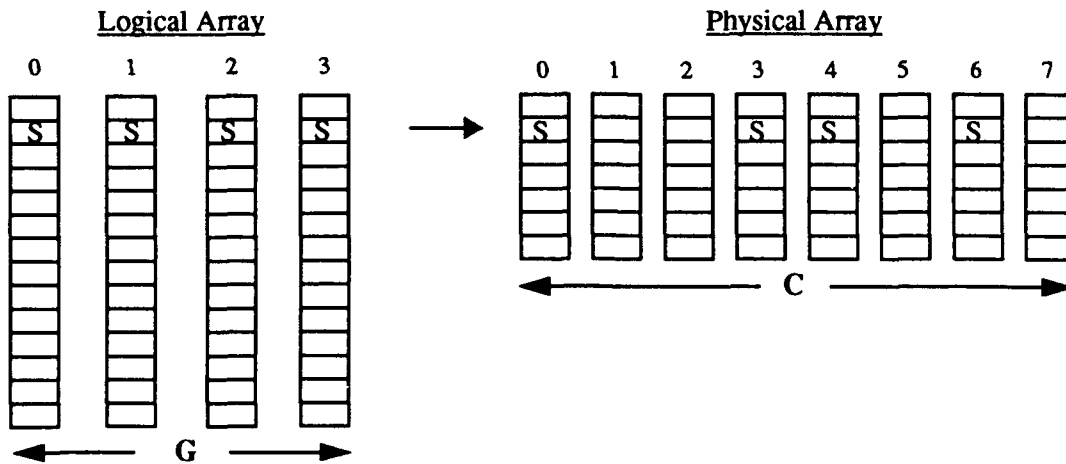
**Figure 2-1: Parity and data layout in a left-symmetric RAID5 organization**

Figure 2-1 shows a layout in which the parity allocation rotates to another disk at each stripe boundary. In general, the parity allocation can be rotated after any number of stripes. It is important to note that a left-symmetric RAID 5 organization defines its data layout to be sequential through parity stripes; that is, user data is logically D0.0, D0.1, D0.2, D0.3, D1.0, D1.1, etc. In general, however, a parity mapping does not imply a data mapping.

In Figure 2-1, parity is computed over the entire width of the array; that is, P0 is the cumulative parity (exclusive-or) of units D0.0 through D0.3. When a disk is identified as failed, any data unit can be reconstructed by reading the corresponding units in the parity stripe, including the parity unit, and computing the cumulative exclusive-or of this data. Note, however, that all the disks in the array are needed by every access that requires reconstruction.

Following Muntz and Lui, let  $G$  be the number of units in a parity stripe, including the parity stripe unit, and consider the problem of decoupling  $G$  from the width of the array. This reduces to a problem of finding a mapping that will allow parity stripes of size  $G$  units to be distributed over some larger number of disks,  $C$ . For our purposes, this larger set of  $C$  disks is the whole array. For comparison purposes, the RAID 5 example in Figure 2-1 has  $G = C = 5$ . This property, that  $G = C$ , defines RAID 5 mappings in the context of this paper.

One perspective on parity declustering in redundant disk arrays is demonstrated in Figure 2-2; a logical RAID 5 array with  $G = 4$  is distributed over  $C=8 > G$  disks, each containing fewer units. The advan-



**Figure 2-2: Declustering a parity stripe of size four over an array of eight disks.**

One advantage of this approach is that it reduces the reconstruction workload applied to each disk during failure recovery. To see this, note that for any given stripe unit on a failed (physical) disk, the parity stripe to which it belongs includes units on only a subset of the total number of disks in the array. In Figure 2-2, for example, disks 1, 2, 5, and 7 do not participate in the reconstruction of the parity stripe marked 'S'. Hence, these disks are called on less often in the reconstruction of one of the other disks. In contrast, a RAID 5 array has  $C = G$ , and so all disks participate in the reconstruction of all lost units of the failed disk.

Figure 2-3 illustrates a declustered parity layout for  $G = 4$  and  $C = 5$ . The procedure for mapping these particular parity stripes to disks is described in Section 4. What is important at this point is that five parity stripes map fifteen data stripe units in the array's first twenty disk units, while in the RAID 5 organization of Figure 2-1, four parity stripes map sixteen data stripe units in the array's first twenty disk units. More disk units are consumed by parity, but not every parity stripe is represented on each disk, so a smaller fraction of each surviving disk is read during reconstruction. For example, if, in Figure 2-3, disk 0 fails, parity stripe four will not have to be read in order to reconstruct it.

Offset	DISK0	DISK1	DISK2	DISK3	DISK4
0	D0.0	D0.1	D0.2	P0	P1
1	D1.0	D1.1	D1.2	D2.2	P2
2	D2.0	D2.1	D3.1	D3.2	P3
3	D3.0	D4.0	D4.1	D4.2	P4

**Figure 2-3: Example data layout in a declustered parity organization.**

Muntz and Lui define the ratio  $(G-1)/(C-1)$  as  $\alpha$ . This parameter, which we refer to as the *declustering*

*ratio*, indicates the degree to which parity stripes are distributed over the array, or, equivalently, it indicates the fraction of each surviving disk that must be read during the reconstruction of a failed disk. Note that  $\alpha = 1$  for the RAID5 organization, indicating that every surviving disk participates in every reconstruction access. All of our performance graphs in Sections 6, 7, and 8 are parameterized by  $\alpha$ .

The parameters  $C$  and  $G$  and the ratio  $\alpha$  together determine the reconstruction performance, the data reliability, and the cost-effectiveness of the array.  $C$  specifies the number of disks in the array, and also determines the data reliability, since larger  $C$  implies a greater chance of a second (data-loss-causing) failure within the array during reconstruction.  $G$ , on the other hand, determines the percentage of total disk space consumed by parity,  $1/G$ . Finally, the declustering ratio  $\alpha$  determines the reconstruction performance of the system; a smaller value should yield better reconstruction performance since a failed disk's reconstruction workload is spread over a larger number of disks. In general, system administrators need to be able to specify  $C$  and  $G$  at installation time according to their cost, performance, capacity, and data reliability needs. This paper provides analyses upon which these decisions can be based.

### 3. Related Work

The idea of improving failure-mode performance by declustering redundancy information originated with mirrored systems [Copeland89, Hsiao90]. Copeland and Keller describe a scheme called *interleaved declustering* which treats *primary* and *secondary* data copies differently. Traditionally, mirrored systems allocate one disk as a primary and another as a secondary. Copeland and Keller instead allocate only half of each disk for primary copies. The other half of each disk contains a portion of the secondary copy data from each of the primaries on all other disks. This insures that a failure can be recovered since the primary and secondary copies of any data are on different disks. It also distributes the workload associated with reconstructing a failed disk across all surviving disks in the array. Hsiao and DeWitt propose a variant of interleaved declustering called *chained declustering*, that increases the array's data reliability.

Muntz and Lui applied ideas similar to those of Copeland and Keller to parity-based arrays. They proposed the declustered parity organization described in Section 2, and then model its reconstruction time analytically, making a number of simplifying assumptions<sup>3</sup>. We attempt in this paper to identify the limits of this theoretical analysis and provide performance predictions based instead on a software implementa-

tion and array simulation. Toward this end we have two primary concerns with the Muntz and Lui analysis.

First, their study assumes that either the set of surviving disks or the replacement disk is driven at 100% utilization. Unfortunately, driving a queueing system such as a magnetic disk at full utilization leads to arbitrarily long response times. Response time is important to all customers and critical in database and on-line transaction-processing (OLTP) systems. An often-cited rule of thumb in the OLTP domain is that 90% of the transactions in a given workload must complete in under two seconds for the system to be acceptable [Anon85, TPCA89]. In a continuous-operation system that requires minutes to hours for the recovery of a failed disk, this rule will apply even during these relatively rare recovery intervals. Our analysis reports on user response time during recovery and presents a simple scheme trading off reconstruction time for user response time.

Our second concern with the Muntz and Lui analysis is that their modeling technique assumes that all disk accesses have the same service time distribution. Unfortunately, real disk accesses are subject to positioning delays that are dependent on the current head position and the position of target data. As an example, suppose that a given track on a replacement disk is being reconstructed, and that a few widely scattered stripe units on that track are already valid because they were written as a result of user (not reconstruction) accesses during reconstruction. These units may either be skipped over by reconstruction, or they may simply be reconstructed along with the rest of the track and over-written with the data that they already hold. To skip over them requires multiple disk accesses instead of one, which is quite likely to cause rotation slips [Chen90b]. In either case, the fact that these sectors have been previously reconstructed is not going to speed the reconstruction process. The Muntz and Lui model assumes that reconstruction time is reduced by a factor equal to the size of the units not needing reconstruction divided by the size of the track, which is not the case. This idea that disk drives are not "work-preserving" due to head positioning and spindle rotation delays is an effect that is difficult to model analytically, but relatively straightforward to address in a simulation-based study.

In this paper we will use balanced incomplete and complete block designs<sup>4</sup> to achieve better perfor-

---

3. Menon and Kasson [Menon92] also proposed a layout where the number of disks exceeds the number of stripe units in a parity stripe, but they mention it only in passing, and do not analyze its benefits and drawbacks.

4. Block designs and their parameters are described in Section 4.2.



mance during reconstruction. Reddy [Reddy91] has also used block designs for improving the degraded-mode performance of a disk array. His organization uses a block design containing  $b$  tuples on  $C$  objects to divide the array into exactly two parity groups: track  $j$  on disk  $i$  is a member of parity group one if object  $i$  is a member of block  $(j \bmod b)$ , where  $b$  is the size of the block design, and is a member of parity group zero otherwise. This generates a layout with properties similar to ours, but is restricted to the case where  $G = C/2$ .

## 4. Data Layout Strategy

### 4.1. Layout Goals

Previous work on declustered parity has left open the problem of allocating parity stripes in an array. Extending from non-declustered parity layout research, we have identified six criteria for a good parity layout [Lee90, Dibble90]. The first four of these deal exclusively with relationships between data stripe units, parity stripe membership, and disk allocation, while the last two of these make recommendations for the relationship between user data allocation and parity stripe organization. Because file systems are free to and often do allocate user data arbitrarily into whatever logical space a storage subsystem presents, our parity layout procedures have no control over these latter two criteria.

1. *Single Failure Correcting.* No two stripe units in the same parity stripe are allowed to reside on the same physical disk. This is the basic characteristic of any redundancy organization that recovers the data of failed disks. In arrays that inter-link disks by a common failure mode, such as power or data cabling, this criteria should be extended to prohibit the allocation of stripe units from one parity stripe to two or more disks sharing that common failure mode [Schulze89].
2. *Distributed Reconstruction.* When any disk fails, its user workload should be evenly distributed across all other disks in the array. When the disk is replaced or repaired, its reconstruction workload should also be evenly distributed.
3. *Distributed Parity.* Parity information should be evenly distributed across the array. All data updates cause a parity update, and so an uneven parity distribution would lead to hot-spot contention, since the disks with more parity would experience more load.
4. *Efficient Mapping.* The functions mapping a file system's logical block address to physical disk address(es) for the corresponding data stripe unit(s) and parity stripe(s) (and the appropriate inverse mappings) must be efficiently implementable; they should consume neither excessive CPU time nor

memory.

5. *Large Write Optimization.* The allocation of contiguous user data to stripe units should correspond to the allocation of stripe units to parity stripes. This insures that whenever a user performs a write that is the size of the data portion of a parity stripe and starts on a parity stripe boundary, it is possible to execute the write without pre-reading the prior contents of any disk data, since the new parity unit depends only on the new data.
6. *Maximal Parallelism.* A read of contiguous user data with size equal to a stripe unit times the number of disks in the array should induce a single stripe unit read on all disks in the array (while requiring alignment only to a stripe unit boundary). This insures that maximum parallelism can be obtained.

As shown in Figure 2-1, the left-symmetric mapping for RAID 5 arrays meets all of these criteria.

## 4.2. Layout Strategy

Our declustered parity layouts are designed to meet our criterion for distributed reconstruction while also lowering the amount of reconstruction work done by each surviving disk. The distributed reconstruction criterion requires that the same number of parity stripe units be read from each surviving disk during the reconstruction of a failed disk. This will be achieved if the number of times that a pair of disks contain stripe units from the same parity stripe is constant across all pairs of disks. Muntz and Lui recognized and suggested that such layouts might be found in the literature for *balanced incomplete block designs* [Hall86]. This paper demonstrates that this can be done and shows how to do it.

A block design is an arrangement of  $v$  distinct objects into  $b$  tuples<sup>5</sup>, each containing  $k$  elements, such that each object appears in exactly  $r$  tuples, and each pair of objects appears in exactly  $\lambda$  tuples. For example, using non-negative integers as objects, a block design with  $b = 5$ ,  $v = 5$ ,  $k = 4$ ,  $r = 4$ , and  $\lambda = 3$  is given by:

Tuple 0: 0, 1, 2, 3	Tuple 3: 0, 2, 3, 4
Tuple 1: 0, 1, 2, 4	Tuple 4: 1, 2, 3, 4
Tuple 2: 0, 1, 3, 4	

Figure 4-1: Sample Complete Block Design

---

5. These tuples are called *blocks* in the block design literature. We avoid this name as it conflicts with the commonly held definition of a block as a contiguous chunk of data.

This example demonstrates a particularly simple form of block design called a *complete block design* which includes all possible combinations of exactly  $k$  distinct elements selected from the set of  $v$  objects. The number of these combinations is  $\binom{v}{k}$ . It is useful to note that only three of  $v$ ,  $k$ ,  $b$ ,  $r$ , and  $\lambda$  are free variables because the following two relations are always true:  $bk = vr$ , and  $r(k-1) = \lambda(v-1)$ . The first of these relations counts the objects in the block design in two ways, and the second counts the pairs in two ways.

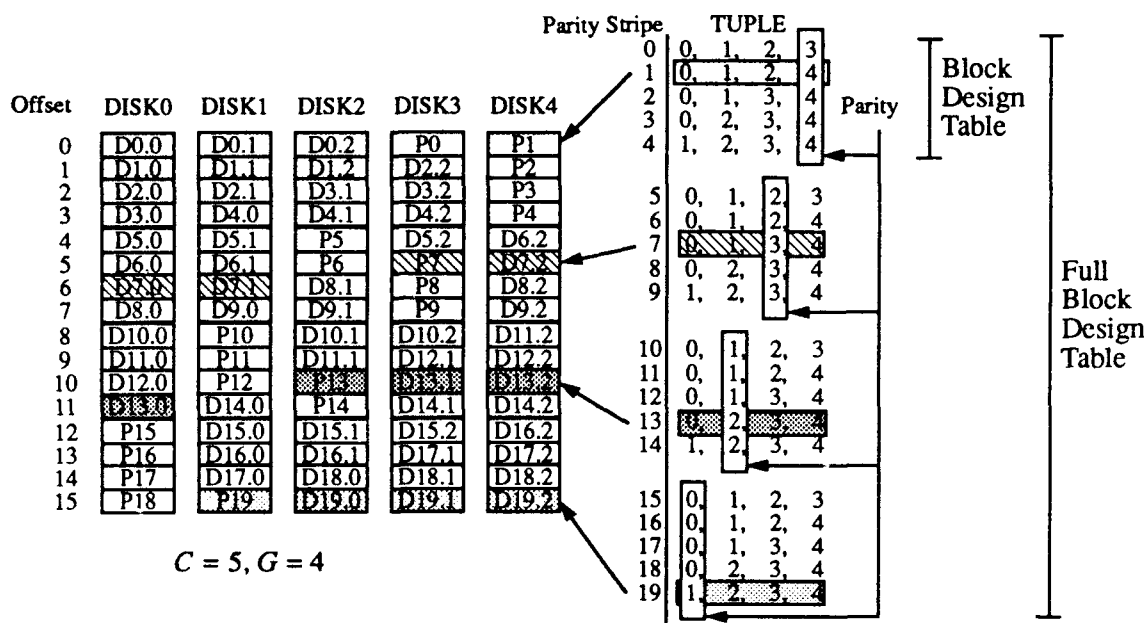
The layout we use associates disks with objects and parity stripes with tuples. For clarity, the following discussion is illustrated by the construction of the layout in Figure 4-2 from the block design in Figure 4-1. To build a parity layout, we find a block design with  $v = C$ ,  $k = G$ , and the minimum possible value for  $b$  (as we shall explain in Section 4.3). Our mapping identifies the elements of a tuple in a block design with the disk numbers on which each successive stripe unit of a parity stripe are allocated. In Figure 4-2, the first tuple in the design of Figure 4-1 is used to lay out parity stripe 0: the three data blocks in parity stripe 0 are on disks 0, 1, and 2, and the parity block is on disk 3. Based on the second tuple, stripe 1 is on disks 0, 1, and 2, with parity on disk 4. In general, stripe unit  $j$  of parity stripe  $i$  is assigned to the lowest available offset on the disk identified by the  $j^{\text{th}}$  element of tuple  $i \bmod b$  in the block design. The layout shown in the top quarter of Figure 4-2, and in Figure 2-3, is derived via this process from the block design in Figure 4-1.

It is apparent from Figure 2-3 that this approach produces a layout that violates our distributed parity criterion (3). To resolve this violation, we derive the layout as above and duplicate it  $G$  times (four times in Figure 4-2), assigning parity to a different element of each tuple in each duplication, as shown in the right side of Figure 4-2. This layout, the entire contents of Figure 4-2, is further duplicated until all stripe units on each disk are mapped to parity stripes. We refer to one iteration of this layout (the first four blocks on each disk in Figure 4-2) as the *block design table*, and one complete cycle (all blocks in Figure 4-2) as the *full block design table*.

We now show how well this layout procedure meets the first four of our layout criteria. Because each tuple contains  $k$  distinct objects, no pair of stripe units in the same parity stripe will be assigned to the same disk. This insures our single failure correcting criterion. The second criterion, which is that reconstruction be distributed evenly, is guaranteed because each pair of objects appears in exactly  $\lambda$  tuples. This means that in each copy of the block design table, disk  $i$  occurs in exactly  $\lambda$  parity stripes with each other disk. Hence, when disk  $i$  fails, every other disk reads exactly  $\lambda$  stripe units while reconstructing the stripe units

## Data Layout on Physical Array

## Layout Derivation from Block Designs



**Figure 4-2: Full block design table for a parity declustering organization.**

associated with each block design table. Note that the actual value of  $\lambda$  is not significant. It is only necessary that it be constant across all pairs of disks, and this is guaranteed by the definition of a block design.

Our distributed parity criterion is achieved because a full block design table is composed of  $G$  block design tables each assigning parity to a different element of each tuple. Referring to Figure 4-2; if we group together the vertical boxes in the right half of the figure we see that the parity assignment function sweeps out the equivalent of one block design table over the course of the full block design table. Since each object appears in exactly  $r$  tuples in a block design table, each will be assigned parity in exactly  $r$  tuples in the full block design table, and so each disk will be assigned  $r$  parity stripe units in every full block design table.

Unfortunately it is not guaranteed that our layout will have an efficient mapping, our fourth criterion, because the size of a block design table is not guaranteed to be small. However, Section 4.3 demonstrates that small block design tables are available for a wide range of parity stripe and array sizes.

Finally, our fifth and sixth criteria depend on the data mapping function used by higher levels of software. Unfortunately, the simple mapping of data to successive data units within successive parity stripes that we use in our simulations, while meeting our large-write optimization criterion, does not meet our

maximal parallelism criterion; that is, not all sets of five adjacent stripe units from the mapping, D0.0, D0.1, D0.2, D1.0, D1.1, D1.2, D2.0, etc., in Figure 4-2 are allocated on five different disks. Instead, reading five adjacent data stripe units starting at stripe unit zero causes disk 0 and 1 to be used twice, and disks 3 and 4 not at all. On the other hand, if we were to employ a data mapping similar to Lee's left-symmetric parity (for non-declustered RAID5 arrays), we may fail to satisfy our large-write optimization criterion. We leave for future work the development of a declustered parity scheme that satisfies both of these criteria.

### 4.3. On the Generation of Block Designs

Complete block designs such as that used in Figure 4-1 are easily generated, but in many cases they are insufficient for our purposes. When the number of disks in an array ( $C$ ) is large relative to the number of stripe units in a parity stripe ( $G$ ) then the size of the block design table becomes unacceptably large and the layout fails our efficient mapping criterion. For example, a 41 disk array with 20% parity overhead ( $G=5$ ) allocated by a complete block design will have about 3,750,000 tuples in its full block design table. In addition to the exorbitant memory requirement for this table, the layout will not meet our distributed parity or distributed reconstruction criteria because even large disks rarely have more than 1,000,000 sectors. For this reason we turn to the theory of *balanced incomplete block designs* [Hall86].

Our goal, then, is to find a small block design on  $C$  objects with a tuple size of  $G$ . This is a difficult problem for general  $C$  and  $G$ . Hall presents a number of techniques, but these are of more theoretical interest than practical value since they do not provide sufficiently general techniques for the direct construction of the necessary designs. Fortunately, Hall also presents a list containing a large number of known block designs, and states that, within the bounds of this list, a solution is given in every case where one is known to exist. Figure 4-3<sup>6</sup> presents a scatter plot of Hall's list of designs. Whenever possible, our parity declustering implementation uses one of Hall's designs.

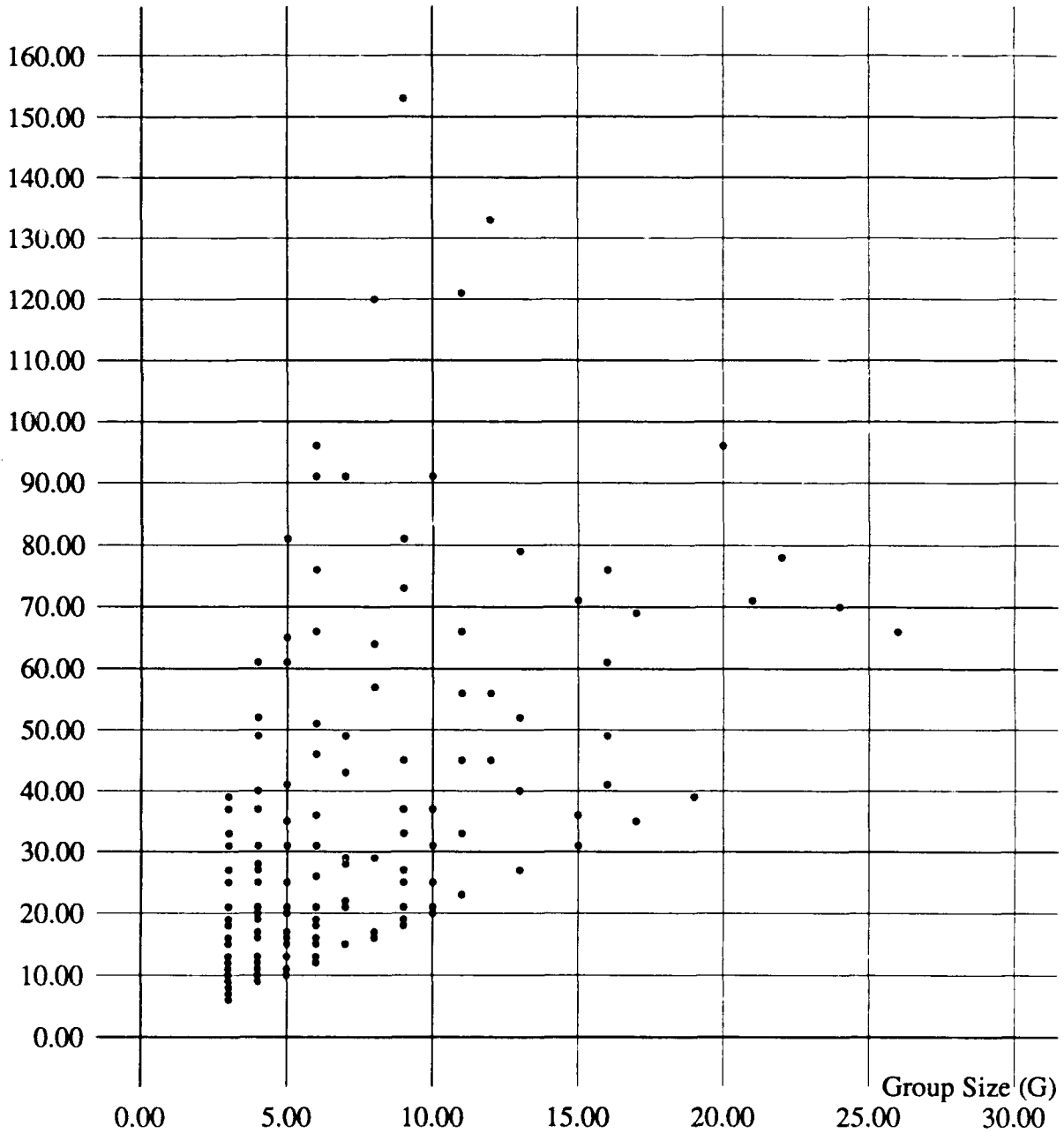
Sometimes a balanced incomplete block design with the required parameters cannot be found. In particular, Hall's table lists only designs with  $G \leq C/2$ , which is essentially equivalent to  $\alpha < 1/2$ . In cases where we cannot find a balanced incomplete block design, we attempt to use a complete block design on

---

6. Note to the reviewer. These (xgraph) figures are not in a satisfactory form for camera-ready submission. Your format suggestions would be much appreciated.

# Figure 4-3: Known Block Designs

Array Size (C)



the indicated parameters. When  $C$  and  $G$  differ by a substantial amount, this method becomes infeasible because the layout table becomes large enough to violate our fourth criterion. In this case, we resort to choosing the closest feasible design point; that is, the point which yields a value of  $\alpha$  closest to what is desired. All of our results indicate that the performance of an array is not highly sensitive to such small variations in  $\alpha$ . The block designs we use in our simulations are given in the appendix.

## 5. The Simulation Environment

We acquired an event-driven disk array simulator called *raidSim* [Chen90b, Lee91] for our analyses. The simulator was developed for the RAID project at U.C. Berkeley [Katz89]. It consists of four primary components. At the top level of abstraction is a synthetic workload generator, which is capable of producing user request streams drawn from a variety of distributions. Table 5-1 (a) shows the configuration of the workload generator used in our simulations. Each request produced by this generator is sent to a *RAID striping driver*, which was originally the actual code used by the Sprite operating system [Ousterhout88] to implement a RAID device on a set of independent disks. Table 5-1 (b) shows the configuration of our extended version of this striping driver. These upper two levels of *raidSim* should actually run on a Sprite machine. Low-level disk operations generated by the striping driver are sent to a *disk simulation module*, which accurately models all significant aspects of each specific disk access (seek time, rotation time, cylinder layout, etc.). Table 5-1 (c) shows the characteristics of the IBM 0661 Model 370 (Lightning) disks on which our simulations are based [IBM0661]. At the lowest level of abstraction in *raidSim* is an *event-driven simulator*, which is invoked to cause simulated time to pass.

It is important to emphasize that the striping driver code was originally taken directly from the Sprite source code, with essentially zero modification to accommodate simulation, and also that our modifications conform to Sprite constraints. This assures that the reference streams generated by the driver are identical to those that would be observed in an actual disk array running the same synthetic workload generator. It also forces us to actually implement our layout strategy and reconstruction optimizations, since we have extended the code in such a way that it could be re-incorporated into an operating system at any time. All reconstruction algorithms discussed in Section 8 have been fully implemented and tested under simulation in our version of the RAID striping driver.

(a) Workload Parameters		(b) Disk Parameters	
Access size:	Fixed at 4 KB	Cylinders:	949
User access rate:	105, 210, and 378 accs/sec	Tracks/Cyl:	14
Alignment:	Fixed at 4 KB	Sectors/Track:	48 @ 512 bytes each
Distribution:	Uniform over all data	Revolution:	13.9 ms
Write Ratio:	0% and 100%(Sections 6, 7)	Seek Time:	2 ms min, 12.5 ms avg, 25 ms max
	50% (Section 8)	Track Skew:	4 sectors

(c) Array Parameters	
Stripe Unit:	Fixed at 4KB
Number of Disks:	Fixed at 21
Head Scheduling:	CVSCAN [Geist87]
Parity Stripe Size:	3, 4, 5, 6, 10, 18, and 21 (RAID5) stripe units
Parity Overhead:	33%, 25%, 20%, 17%, 10%, 6%, and 5%, respectively
Data Layout:	RAID5: Left Symmetric
	Declustered: By parity stripe index
Parity Layout:	RAID5: Left Symmetric
	Declustered: Block Design Based
Power/Cabling:	Each disk independently powered and cabled

**Table 5-1: Simulation Parameters**

## 6. Fault-Free Performance

Figures 6-1 and 6-2 show the average response time experienced by read and write requests in a fault-free disk array as a function of the declustering ratio,  $\alpha$ . Our simulated system has 21 disks, so the fraction of space consumed by parity units,  $1/G$ , is  $1/(20\alpha+1)$ . In the 100% read case we show three average response time curves corresponding to user access rates ("Rate") of 105, 210, and 378 random reads of 4 KB per second (on average, 5, 10, and 18 user reads of 4 KB per second per disk are applied to disks capable of a maximum of about 46 random 4 KB reads per second). In the 100% write case we show two much slower average response time curves corresponding to Rate = 105 and 210 random user writes of 4 KB per second. User writes are much slower than user reads because writes must update parity units as well as data units. Because the RAID striping driver does not have access to higher level software caching, it is unable to exploit caching of the old contents of either data or parity units [Patterson88], and, because it does not have control over the precise timing of its disks, it is unable to execute a parity read-modify-write sequence in one disk access [Menon89]. Without these optimizations, our striping driver's fault-free behavior is to



execute four separate disk accesses for each user write instead of the single disk access needed by an user read. Because of this high cost for user writes, our system is not able to sustain 378 user writes of 4 KB per second (this would be 72 4 KB accesses per second per disk).

Figures 6-1 and 6-2 show that except for writes with  $\alpha = 0.1$ , fault-free performance is essentially independent of parity declustering. This exception is the result of an optimization that the RAID striping driver employs when requests for one stripe unit are applied to parity stripes containing only three stripe units [Chen90a]. In this case the driver can choose to write the specified data unit, read the other data unit in this parity stripe, then directly write the corresponding parity unit in three disk accesses, instead of pre-reading and overwriting both the specified data unit and corresponding parity unit in four disk accesses. In the rest of this paper we will neglect the case of  $\alpha = 0.1$  ( $G = 3$ ) to avoid repeating this optimization discussion.

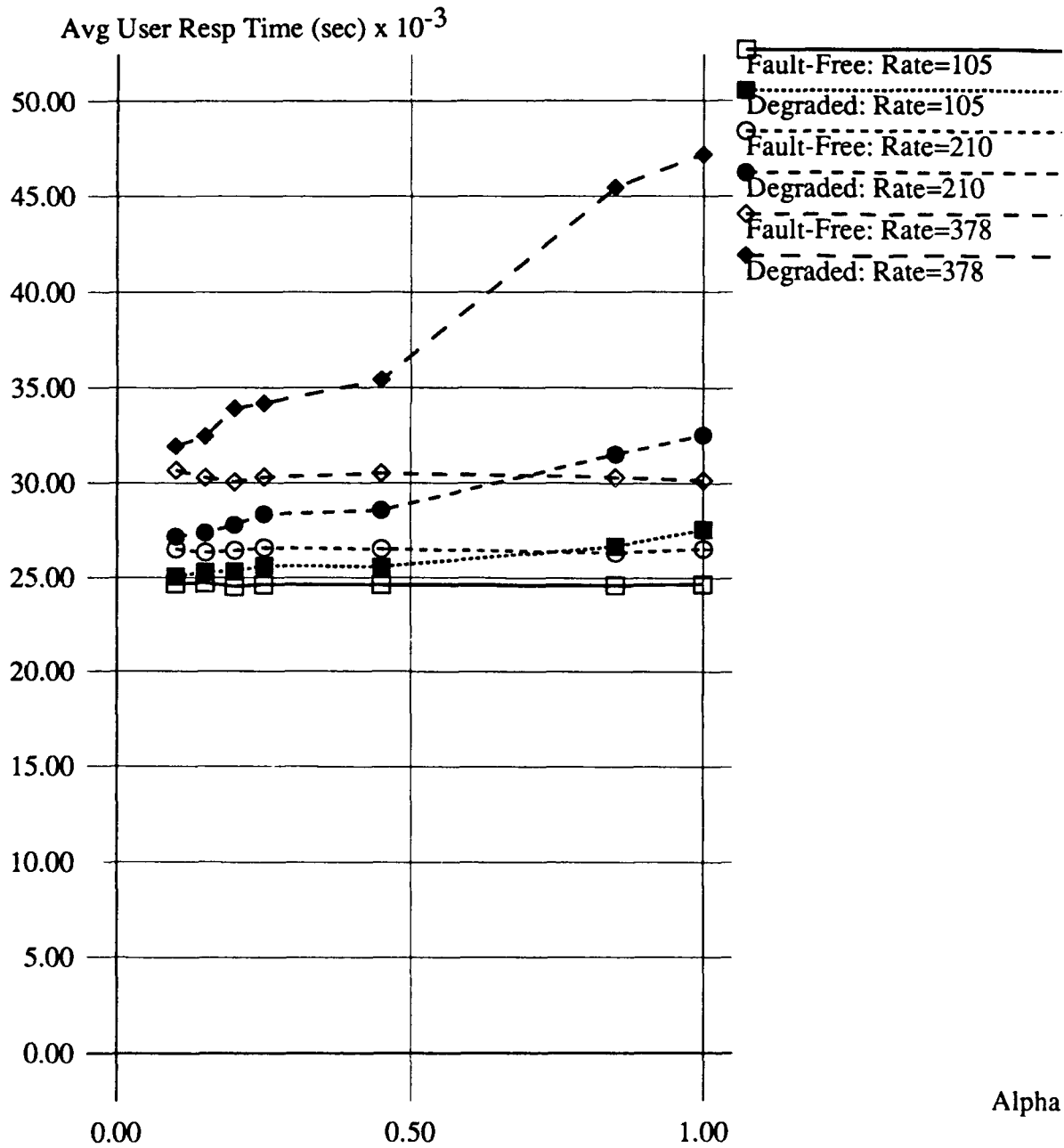
We note that our parity declustering implementation may not perform equivalently to a left-symmetric RAID 5 mapping when user requests are larger than one stripe unit. Declustered parity has the advantage of exploiting our large-write optimization with smaller user writes because it has smaller parity stripes. On the other hand, because our implementation does not currently meet our maximal parallelism criterion, it will not be as able to exploit the full parallelism of the array on large user accesses. Overall performance will be dictated by the balancing of these two effects, and will depend on the access size distribution.

## 7. Degraded-Mode Performance

Figures 6-1 and 6-2 also show average response time curves for our array when it is degraded by a failed disk that has not yet been replaced. In this case, an on-the-fly reconstruction takes place on every access that requires data from the failed disk. When such an access is a read (or pre-read), the failed disk's contents are reconstructed by reading and computing an exclusive-or over all surviving units of the requested data's parity stripe. Because the amount of work this entails depends on the size of each parity stripe, these degraded-mode average response time curves suffer less degradation with a lower parity declustering ratio (smaller  $\alpha$ ). This is the only effect shown in Figure 6-1 (100% Reads), but for writes there is another consideration. When a user write specifies data on a surviving disk whose associated parity unit is on the failed disk, there is no value in trying to update this lost parity. In this case a user write

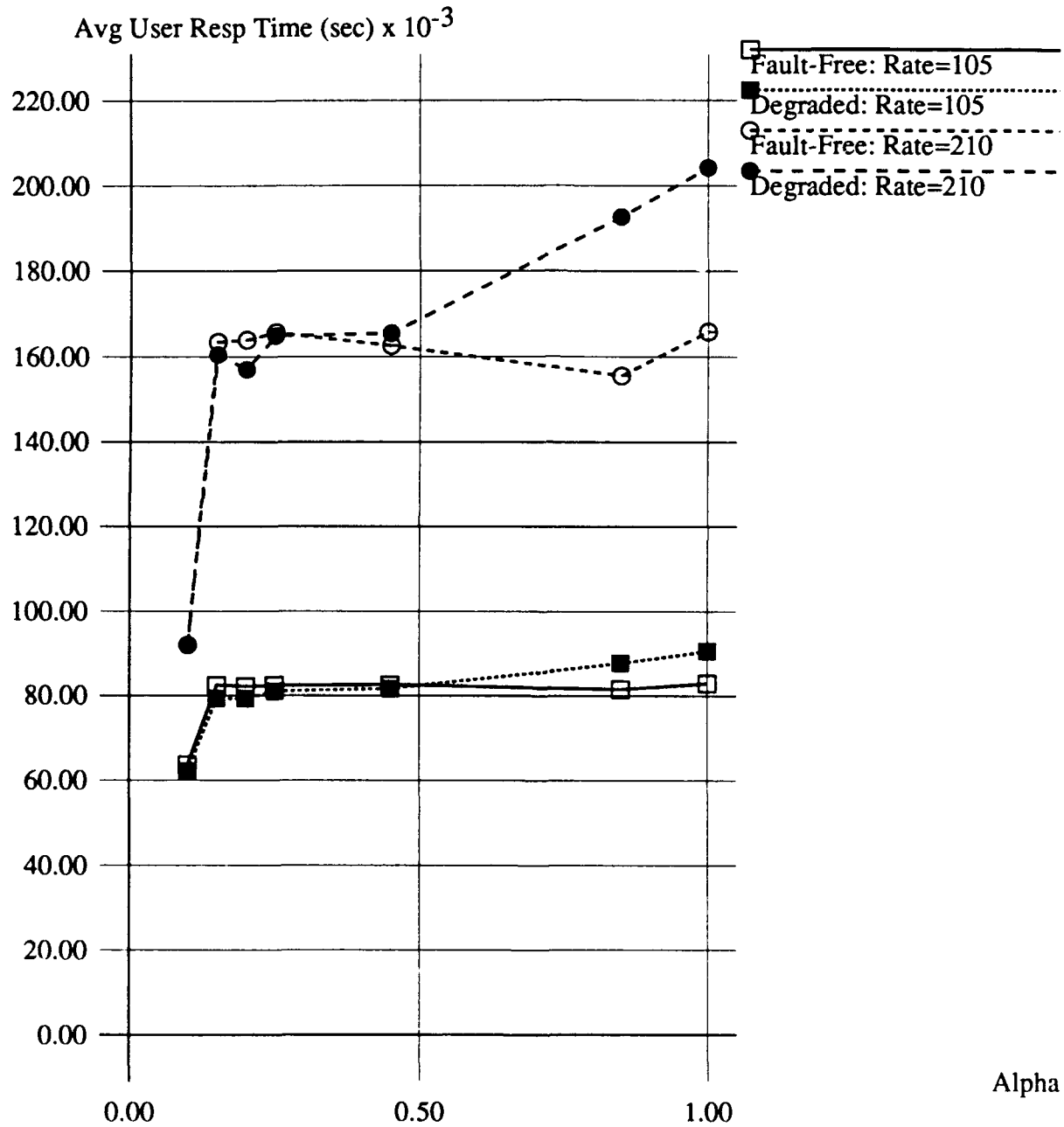
# Figure 6-1: Response Time

100% Reads



# Figure 6-2: Response Time

100% Writes



induces only one, rather than four, disk accesses. This effect decreases the workload of the array relative to the fault-free case, and, at low declustering ratios, it can lead to slightly better average response time in the degraded rather than fault-free mode!

## 8. Reconstruction Performance

The primary purpose for parity declustering over RAID 5 organizations has been given as the desire to support higher user performance during recovery and shorter recovery periods [Muntz90]. In this section we show this to be effective, although we also show that previously proposed optimizations are not entirely so. We also demonstrate that there remains an important trade-off between higher performance during recovery and shorter recovery periods.

Reconstruction, in its simplest form, involves a single sweep through the contents of a failed disk. For each stripe unit on a replacement disk, the reconstruction process reads all other stripe units in the corresponding parity stripe and computes an exclusive-or on these units. The resulting unit is then written to the replacement disk. The time needed to entirely recover a failed disk is equal to the time needed to replace it in the array plus the time needed to reconstruct its entire contents and store them on the replacement. This latter time is termed the *reconstruction time*. In an array that maintains a pool of on-line spare disks, the replacement time can be kept sufficiently small that repair time is essentially reconstruction time. Highly available disk arrays require short repair times to assure high data reliability because the mean time until data loss is inversely proportional to mean repair time [Patterson88]. However, minimal reconstruction time occurs when user access is denied during reconstruction. Because this cannot take less than the three minutes it takes to read all sectors on our disks into their track buffers, and usually takes much longer, continuous-operation systems require data availability during reconstruction.

Muntz and Lui identify two optimizations to a simple sweep reconstruction: *redirection of reads* and *piggybacking of writes*. In the first, user accesses<sup>7</sup> to data that has already been reconstructed are serviced by (redirected to) the replacement disk, rather than invoking on-the-fly reconstruction as they would if the data were not yet available. This reduces the number of disk accesses necessary to service a typical request

---

7. In what follows, we distinguish between *user accesses*, which are those generated by applications as part of the normal workload, and *reconstruction accesses*, which are those generated by a background reconstruction process to regenerate lost data and store it on a replacement disk.

during reconstruction. In the second optimization, user read-accesses that cause on-the-fly reconstruction also cause the reconstructed data to be written to the replacement disk. This is targeted at speeding reconstruction, since those units need not be subsequently reconstructed.

Muntz and Lui also mention that in servicing a user's write to a stripe unit whose contents have not yet been reconstructed, the device driver has a choice of writing the new data directly to the replacement disk or performing an on-the-fly reconstruct-read followed by an update of the associated parity unit alone. Their model assumes that users' writes are always sent to the replacement disk if appropriate, but for reasons explained below, we question whether this is always a good idea. Therefore, we investigate four algorithms, distinguished by the amount and type of non-reconstruction workload they send to the replacement disk. In our *baseline algorithm*, no extra work is sent; whenever possible user writes are folded into the parity unit, and neither reconstruction optimization is enabled. In our *user-writes* algorithm, only user writes explicitly targeted at the replacement disk are sent directly to the replacement. The *redirection* algorithm adds the redirection of reads optimization to the user-writes case. Finally, the *redirect plus piggyback* algorithm adds the piggybacking of writes optimization to the redirection algorithm.<sup>8</sup>

## 8.1 Single Threaded vs. Parallel Reconstruction

Figures 8-1 and 8-2 present the reconstruction time and average user response time for our four reconstruction algorithms under a user workload that is 50% 4 KB random reads and 50% 4 KB random writes. These figures show the substantial effectiveness of parity declustering for lowering both the reconstruction time and average user response time relative to a RAID 5 organization ( $\alpha = 1.0$ ). For example, at 105 user accesses per second, an array with declustering ratio 0.15 reconstructs a failed disk about twice as fast as the RAID 5 array while delivering an average user response time that is about 33% lower.

While Figures 8-1 and 8-2 make a strong case for the advantages of declustering parity, even the fastest reconstruction shown, about 60 minutes<sup>9</sup>, is 20 times longer than the physical minimum time these disks take to read or write their entire contents. The problem here is that a single reconstruction process is not able to highly utilize any of the array's disks, particularly at low declustering ratios. To further reduce

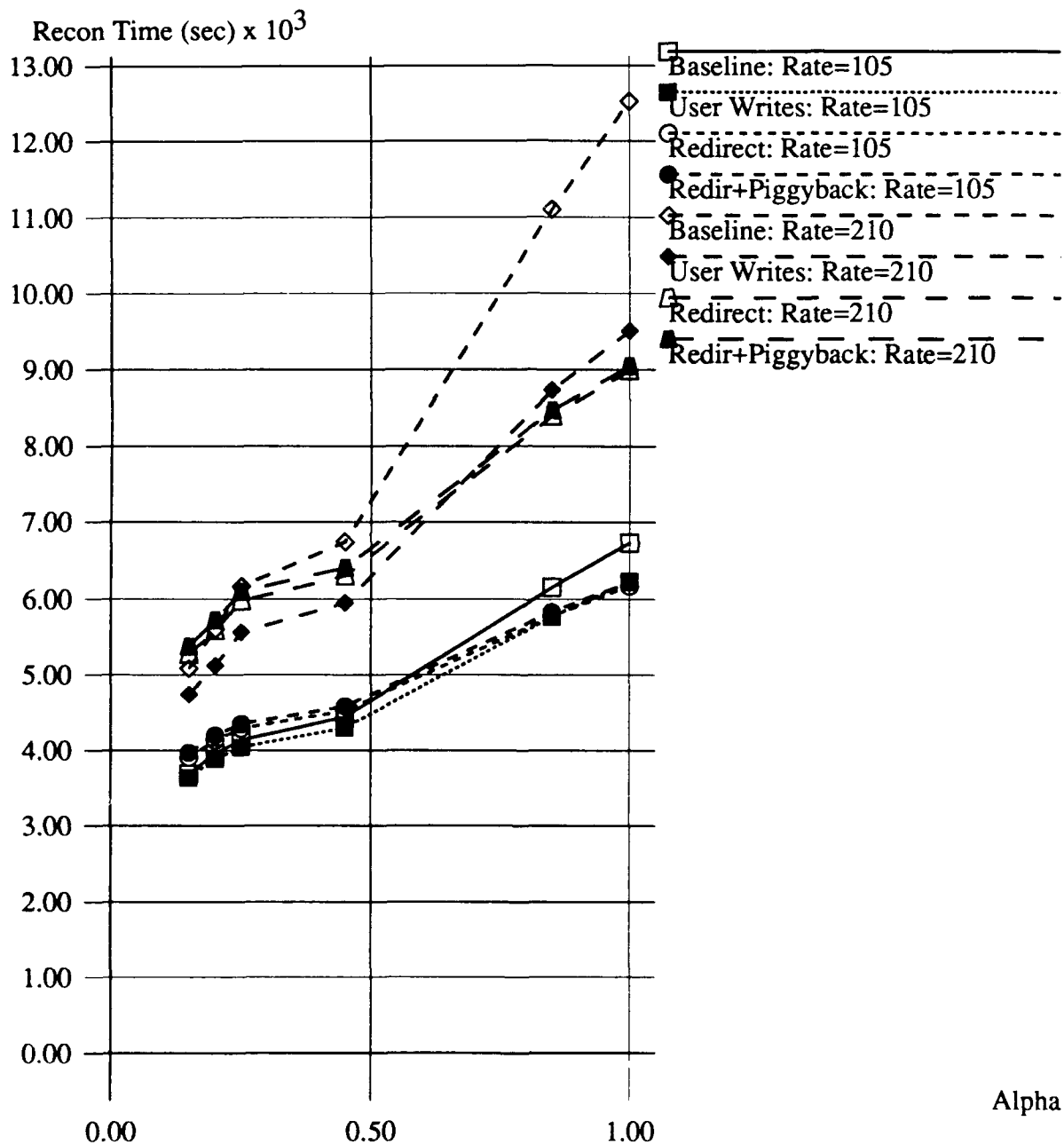
---

8. Note that our definition of a 'baseline' differs from that of Muntz and Lui, who use the term to identify what we call the user-writes case.

9. For comparison, RAID products available today specify on-line reconstruction time to be in the range of one to four hours [Meador89, Rudeseal92].

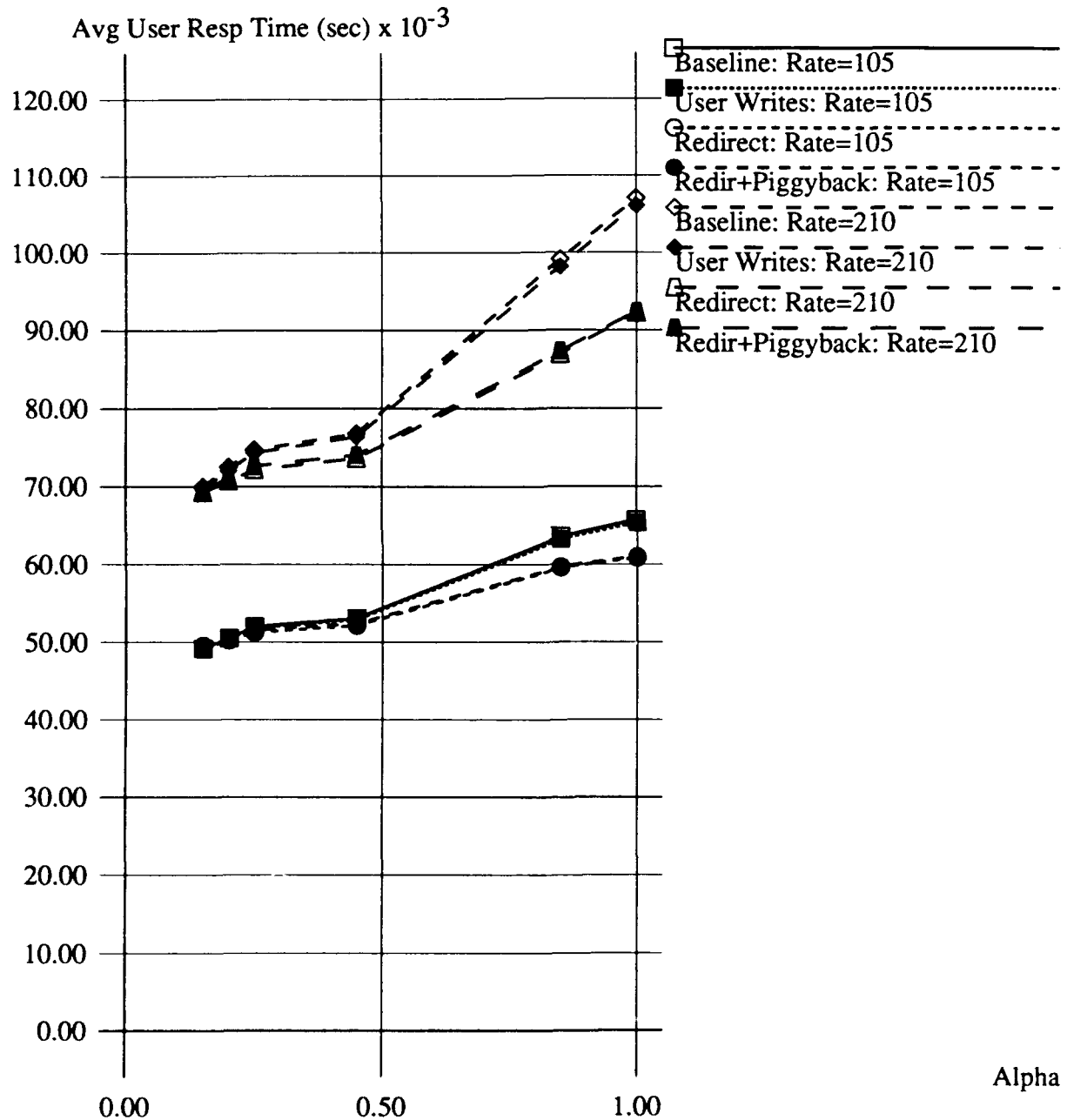
# Figure 8-1: Single-Thread Reconstruction Time

50% Reads, 50% Writes



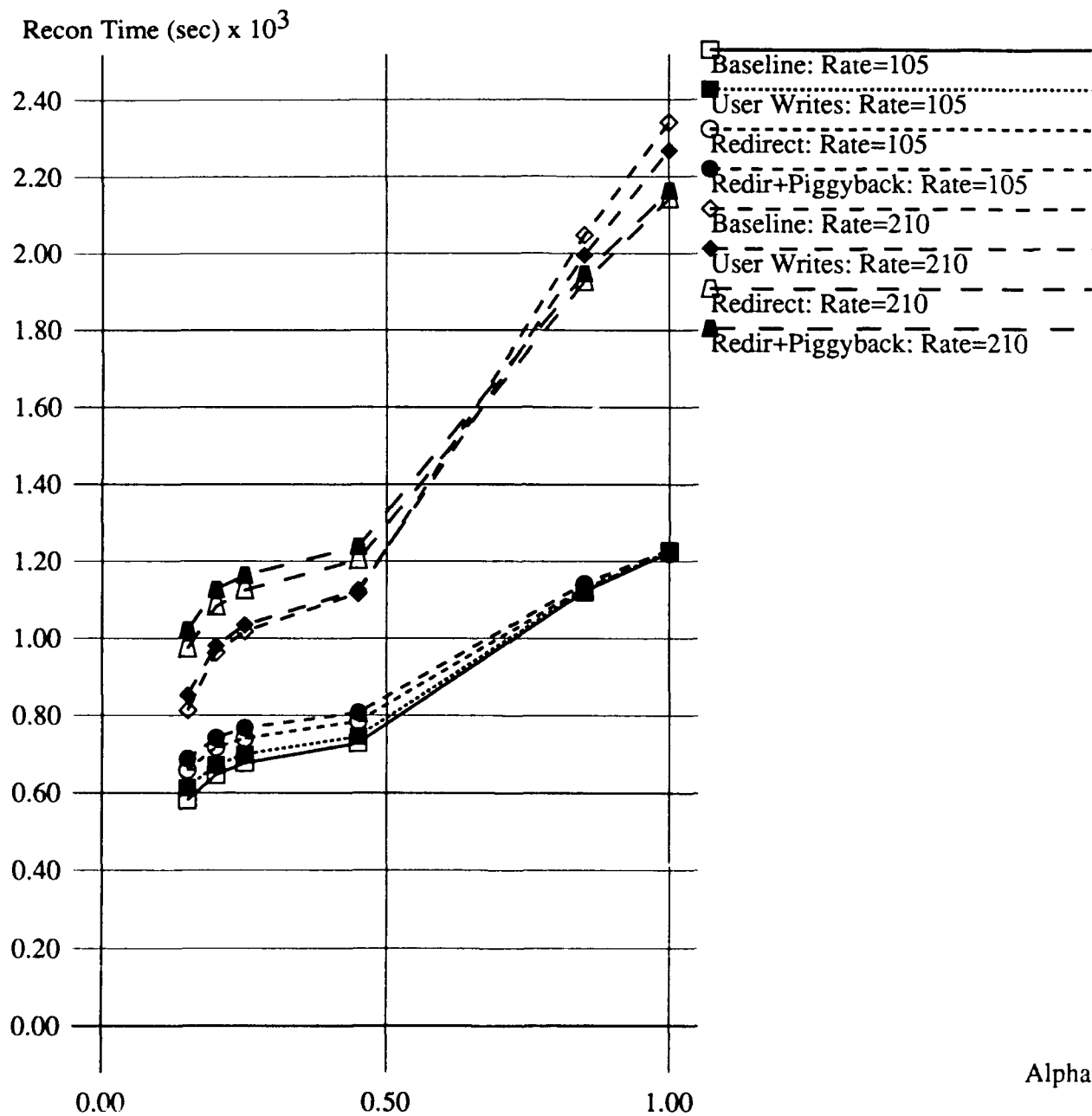
# Figure 8-2: Single-Thread Avg User Response Time

50% Reads, 50% Writes



# Figure 8-3: Eight-Way Parallel Reconstruction Time

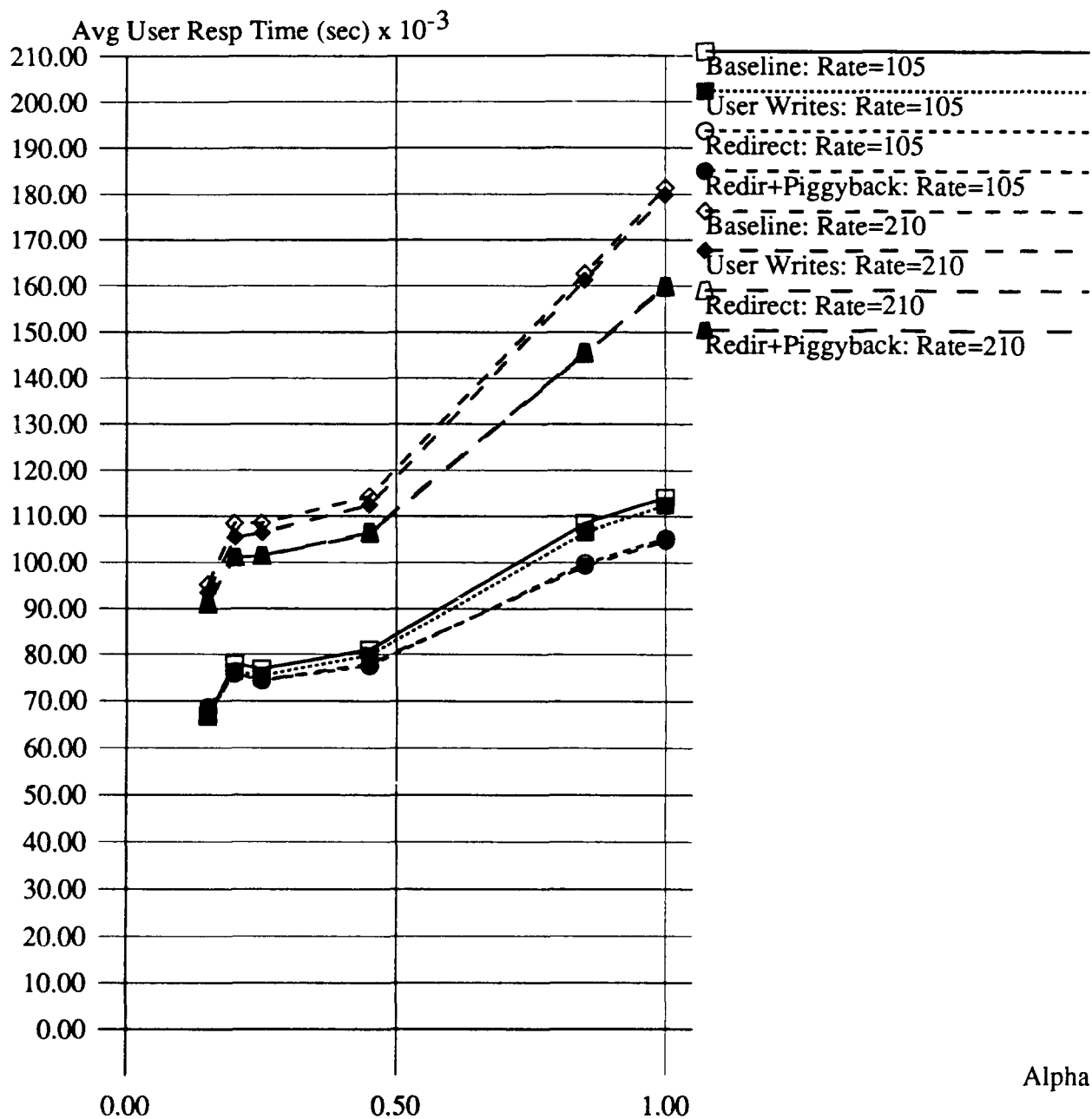
50% Reads, 50% Writes





# Figure 8-4: Eight-Way Parallel Avg User Response Time

50% Reads, 50% Writes



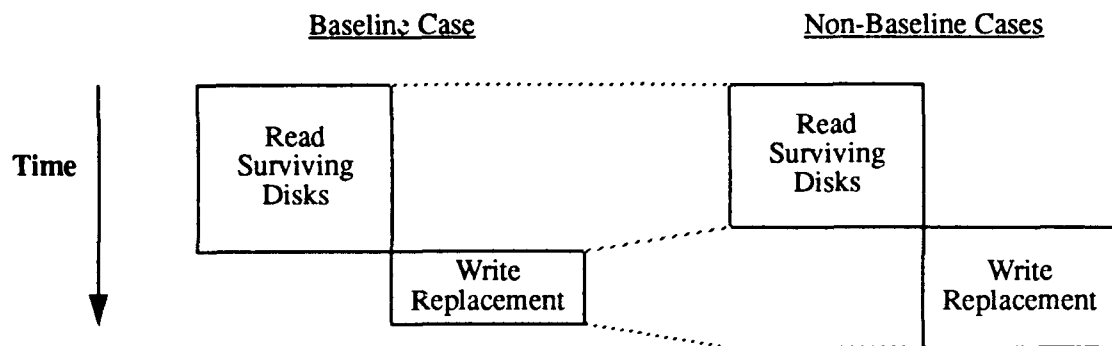
reconstruction time we introduce parallel reconstruction processes. This should speed reconstruction substantially, although it will also further degrade the response time of concurrent user accesses.

Figures 8-3 and 8-4 show the reconstruction time and average user response time when eight processes are concurrently reconstructing the replacement disk. For these workloads of 105 and 210 user accesses per second, reconstruction time is reduced by a factor of four to six relative to Figure 8-1. This gives reconstruction times between 10 and 40 minutes. However, response times suffer increases of 35% to 75%. Still, even the worst of these average response times is less than 200 ms, so a simple transactions such as TPCA, which should require less than three disk accesses per transaction, should have a good chance of meeting its required transaction response time of two seconds.

## 8.2 Comparing Reconstruction Algorithms

The response time curves of Figures 8-2 and 8-4 show that Muntz and Lui's redirection of reads optimization has little benefit in lightly-loaded arrays with a low declustering ratio, but can benefit heavily-loaded RAID 5 arrays with 10% to 15% reductions in response time. The addition of piggybacking of writes to the redirection of reads algorithm is intended to reduce reconstruction time without penalty to average user response time. With the user loads we employ, piggybacking of writes yields very little improvement of penalty over redirection of reads alone. We will not pursue piggybacking of writes further in this section.

Figures 8-1 and 8-2 also show that the two "more optimized" reconstruction algorithms do not consistently decrease reconstruction time relative to the simpler algorithms. In particular, the single-threaded user-writes algorithm yields faster reconstruction times than all others for all values of  $\alpha$  less than 0.5. Similarly, the eight-way parallel baseline and user-writes algorithms yields faster reconstruction times than the other two for all values of  $\alpha$  less than 0.5. These are surprising results. We had expected the more optimized algorithms to experience improved reconstruction time due to the off-loading of work from the over-utilized surviving disks to the under-utilized replacement disk, and also because in the piggybacking case, they reduce the number of units that need to get reconstructed. The reason for this reversal is that loading the replacement disk with random work penalizes the reconstruction writes to this disk more than off-loading benefits the surviving disks unless the surviving disks are highly utilized.



**Figure 8-5: The effect of the optimizations on the reconstruction cycle.**

Because reconstruction time is the time taken to reconstruct all parity stripes associated with the failed disk's data, one at a time, our unexpected effect can be understood by examining the number of and time taken by reconstructions of a single unit as shown in Figure 8-5. We call this single stripe unit reconstruction period a *reconstruction cycle*. It is composed of a read phase and a write phase. The length of the read phase, the time needed to collect and exclusive-or the surviving stripe units, is approximately the maximum of  $G-1$  reads of disks that are also supporting a substantial load of random user requests. If a small amount of work is off-loaded from these disks while they are not excessively loaded, it will not significantly reduce this maximum access time. But, even a small amount of random load imposed on the replacement disk may greatly increase its average access times because reconstruction writes are sequential and do not require long seeks. This effect, suggesting a preference for algorithms that minimize non-reconstruction activity in the replacement disk, must be contrasted with the reduction in number of reconstruction cycles that occurs when user activity causes writes to the portion of the replacement disk's data that has not yet been reconstructed.

Table 8-1 presents a sample of the durations of the intervals in Figure 8-5. These numbers are averaged over the reconstruction of the last 300 stripe units on a replacement disk. The standard deviations are shown in parentheses. For these final reconstruction cycles, the piggybacking of writes is not likely to occur, but the redirection of reads will be at its greatest utility. This table shows that the more complex algorithms tend to yield lower read phase times and higher write phase times. These numbers suggest that with a low declustering ratio the baseline algorithm has an advantage over the user-writes algorithm and both are faster than the other two algorithms.

This suggestion is not entirely borne out by Figures 8-1 and 8-3 because the baseline algorithm gets

none of its reconstruction work done for it by user requests, as is the case for the other algorithms. In the single threaded case, reconstruction is so slow that this latter effect dominates and the baseline algorithm reconstructs more slowly than the others for all but the smallest declustering ratio. However, in the eight-way parallel case, reconstruction is fast enough that this “free reconstruction” effect does not compensate for longer services times experienced by the replacement disk because free reconstruction moves the heads around randomly.

Single-Thread Reconstruction			
	$\alpha = 0.15$	$\alpha = 0.45$	$\alpha = 1.0$
Baseline	$88(2)+15(0.2) = 103$	$125(4)+15(0.2) = 140$	$198(8)+15(0.1) = 213$
User-Writes	$68(2)+17(0.02) = 85$	$97(2)+19(0.02) = 116$	$196(5)+23(0) = 219$
Redirect	$65(2)+50(2) = 115$	$86(2)+41(0.1) = 127$	$137(3)+42(1) = 179$
Redir+Piggyback	$64(3)+44(2) = 109$	$86(2)+4(2) = 127$	$134(4)+46(0.1) = 180$

Eight-Way Parallel Reconstruction			
	$\alpha = 0.15$	$\alpha = 0.45$	$\alpha = 1.0$
Baseline	$69(4)+10(0.4) = 79$	$96(4)+11(0.2) = 107$	$204(8)+14(0) = 218$
User-Writes	$89(3)+27(0.4) = 116$	$121(3)+23(0.2) = 144$	$225(8)+22(0) = 247$
Redirect	$85(3)+58(2) = 143$	$102(2)+50(3) = 152$	$160(4)+36(0.5) = 196$
Redir+Piggyback	$85(1)+58(2) = 143$	$102(2)+50(0.9) = 152$	$160(2)+36(1) = 196$

**Table 8-1: Reconstruction cycle times (ms) at Rate = 210**  
*read\_time(std. dev) + write\_time(std. dev.) = cycle\_time.*

### 8.3 Comparison to Analytic Model

Muntz and Lui also proposed an analytic expression for reconstruction time in an array employing declustered parity [Muntz90]. Figure 8-6 shows our best attempt to reconcile their model with our simulations. Because their model takes as input the fault-free arrival rate and read-fraction of accesses to the disk rather than of user requests, we apply the following conversions, required because each user write induces two disk reads and two disk writes. Letting the fraction of user accesses that are reads be  $R$ , Muntz and Lui’s rate of arrival of disk accesses is  $(4-3R)$  times larger than our user request arrival rate, and the fraction of their disk accesses that are reads is  $(2-R)/(4-3R)$ .

Their model also requires as input the maximum rate at which each disk executes its accesses. It is this parameter that causes most of the disagreement between their model and our simulations. In Figure 8-6 we

use the maximum rate at which one disk services entirely random accesses (46 accesses per second) in their model. This is consistent with the way Muntz and Lui's designed their model, but it does not account for the fact that real disks execute sequential accesses much faster than random accesses. For example, with these disks and random accesses the minimum time required to read or write an entire disk is over 1700 seconds - more than three times longer than our fastest simulated reconstruction time.

The other manifestation of the problem of using a single disk service rate is their predictions for the benefits of the redirection of reads and piggybacking of writes optimizations. Because redirecting work to the replacement disk in their model does not increase this disk's average access time (as is shown to be false for our simulations in Table 8-1), their predictions for the user-writes algorithm are more pessimistic than for their other algorithms.

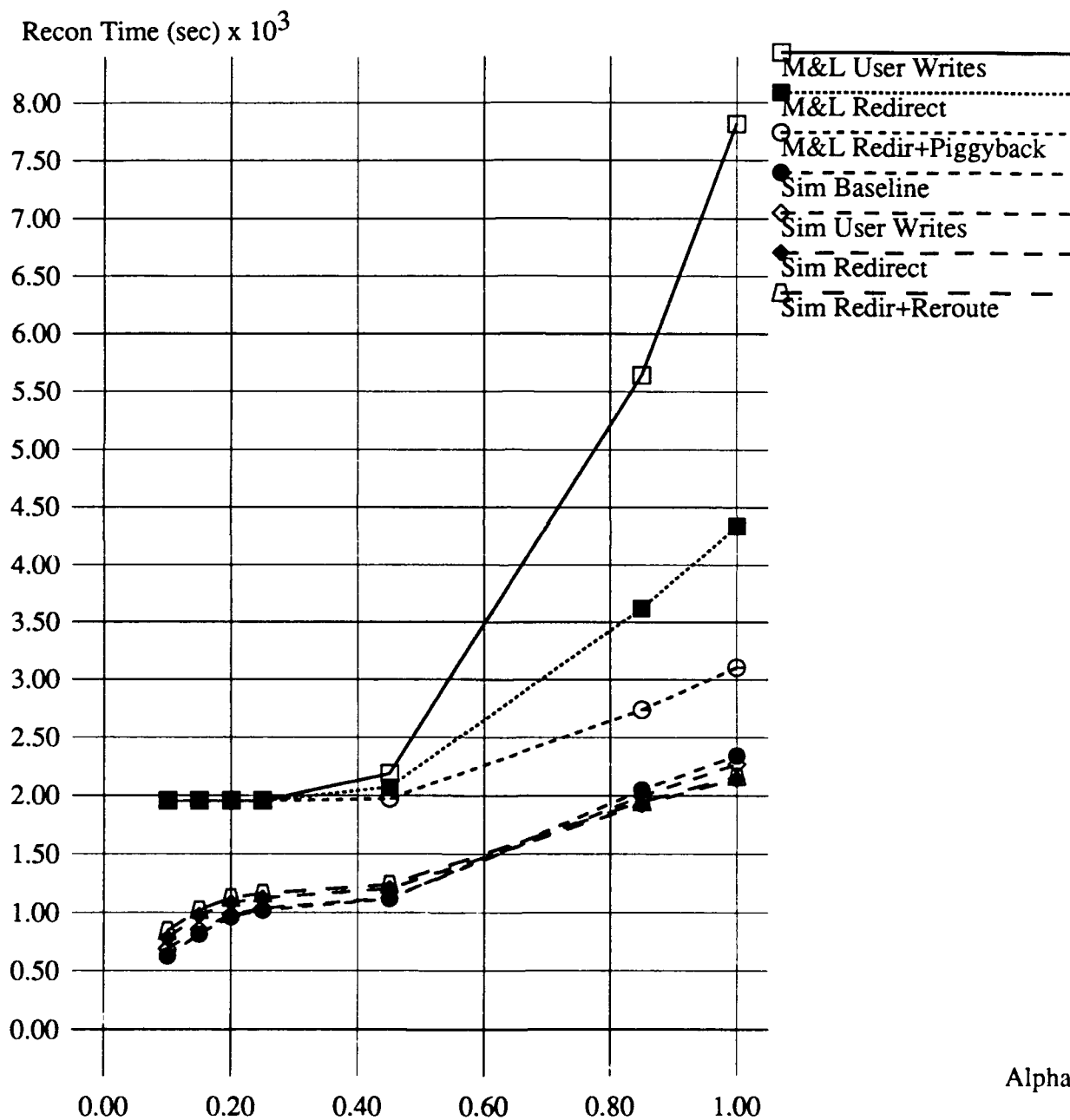
## 9. Conclusions and Future Work

In this paper we have demonstrated that parity declustering, a strategy for allocating parity in a single-failure-correcting redundant disk array that trades increased parity overhead for reduced user performance degradation during on-line failure recovery, can be effectively implemented in array controlling software. We have exploited the special characteristics of balanced incomplete and complete block designs to provide array configuration flexibility while meeting most of our proposed criteria for the "goodness" of a parity allocation. In particular, using block designs to map parity stripes onto a disk array insures that both the parity update load and the on-line reconstruction load is balanced over all (surviving) disks in the array. This is achieved without performance penalty in a normal (non-recovering) user workload dominated by small (4 KB) disk accesses.

Previous work proposing declustered parity mappings for redundant disk arrays has suggested without implementation the use of block designs, has proposed two optimizations to a simple sweep reconstruction algorithm, and has analytically modeled the resultant expected reconstruction time [Muntz90]. In addition to extending their proposal to an implementation, our work evaluates their optimizations and reconstruction time models in the context of our software implementation running on a disk-accurate simulator. Our findings are strongly in agreement with theirs about the overall utility of parity declustering, but we disagree with their projections for expected reconstruction time and the value of their optimized recon-

# Figure 8-6: Comparing M&L Model to Simulation

210 User Access/Second



struction algorithms. Our disagreements largely result because of our more accurate models of magnetic-disk accesses. We find that their estimates for reconstruction time are significantly pessimistic and that their optimizations actually slow reconstruction in some cases. We also find that in practice, multiple, parallel reconstruction processes are necessary to attain fast reconstruction, although this additional load can degrade user response time substantially. Our most surprising result is that in an array with a low declustering ratio and parallel reconstruction processes, the simplest reconstruction algorithm produces the fastest reconstruction time because it best avoids disk seeks.

We feel that future work on parity declustering and related performance/reliability trade-offs is rich in research problems, a sample of which we give here. Our use of block designs would be greatly improved if we could find a wider range of parameters. In particular, satisfactorily small block designs with a declustering ratio between 0.5 and 0.8 are unknown to us. We have also not spent much time on the data mapping in our declustered parity arrays; we think many of our layouts may be amenable to both our large write optimization and maximal parallelism criteria. With an eye to our performance modeling, we would like to see Muntz and Lui's analytical model modified to incorporate more of the complexity of disk accesses and we intend to explore disk arrays with different stripe unit sizes and user workload characteristics. One important concern that neither our nor Muntz and Lui's work considers is the impact of CPU overhead and architectural bottlenecks in the reconstructing system [Chervenak91]. For greater control of the reconstruction process we intend to implement throttling of reconstruction and/or user workload as well as a flexible prioritization scheme that reduces user response time degradation without starving reconstruction. Finally, we hope to install our software implementation of parity declustering on an experimental, high-performance redundant disk array and measure its performance directly.

## References

- [Anon85] Anon, et. al., "A Measure of Transaction Processing Power," *Datamation*, Vol. 31.7, April 1985.
- [Bitton88] D. Bitton and J. Gray, "Disk Shadowing," *Proceedings of the 14th Conference on Very Large Data Bases*, 1988.
- [Chen90a] P. Chen, et. al., "An Evaluation of Redundant Arrays of Disks using an Amdahl 5890," *Proceedings of the ACM SIGMETRICS Conference*, 1990.
- [Chen90b] P. Chen and D. Patterson, "Maximizing Performance in a Striped Disk Array", *Proceedings of ACM SIGARCH Conference*, 1990.

- [Chervenak91] A. Chervenak and R. Katz, "Performance of a RAID Prototype," *Proceedings of the ACM SIGMETRICS Conference*, 1991.
- [Copeland89] G. Copeland and T. Keller, "A Comparison of High-Availability Media Recovery Techniques," *Proceedings of the ACM SIGMOD Conference*, 1989.
- [Dibble90] P. Dibble, "A Parallel Interleaved File System," University of Rochester, Technical Report 334, 1990.
- [Geist87] R. Geist and S. Daniel, "A Continuum of Disk Scheduling Algorithms," *ACM Transactions on Computer Systems*. Vol. 5(1), 1987.
- [Gray90] G. Gray, B. Horst, and M. Walker, "Parity Striping of Disc Arrays: Low-Cost Reliable Storage with Acceptable Throughput," *Proceedings of the 16th Conference on Very Large Data Bases*, 1990.
- [Hall86] M. Hall, *Combinatorial Theory*, Wiley-Interscience, 1986.
- [Hsiao90] H.-I. Hsiao and D. DeWitt, "Declustering: A New Availability Strategy for MultiProcessor Database Machines," *Proceedings of the 6th International Data Engineering Conference*, 1990.
- [IBM0661] IBM Corporation, IBM 0661 Disk Drive Product Description, Model 370, First Edition, Low End Storage Products, 504/114-2, 1989.
- [Katz89] R. Katz, et. al., "A Project on High Performance I/O Subsystems," *ACM Computer Architecture News*, Vol. 17(5), 1989.
- [Kim86] M. Kim, "Synchronized Disk Interleaving," *IEEE Transactions on Computers*, Vol. 35 (11), 1986.
- [Lee90] E. Lee, "Software and Performance Issues in the Implementation of a RAID Prototype," University of California, Technical Report UCB/CSD 90/573, 1990.
- [Lee91] E. Lee and R. Katz, "Performance Consequences of Parity Placement in Disk Arrays," *Proceedings of ASPLOS-IV*, 1991.
- [Livny87] M. Livny, S. Khoshafian, and H. Boral, "Multi-Disk Management Algorithms," *Proceedings of the ACM SIGMETRICS Conference*, 1987.
- [Meador89] W. Meador, "Disk Array Systems," *Proceedings of COMPCON*, 1989.
- [Menon92] J. Menon and J. Kasson, "Methods for Improved Update Performance of Disk Arrays," *Proceedings of the Hawaii International Conference on System Sciences*, 1992.
- [Muntz90] R. Muntz and J. Lui, "Performance Analysis of Disk Arrays Under Failure," *Proceedings of the 16th Conference on Very Large Data Bases*, 1990.
- [Ousterhout88] J. Ousterhout, et. al., "The Sprite Network Operating System," *IEEE Computer*, February 1988.
- [Patterson88] D. Patterson, G. Gibson, and R. Katz, "A Case for Redundant Arrays of Inexpensive Disks (RAID)," *Proceedings of the ACM SIGMOD Conference*, 1988.



[Reddy89] A. Reddy and P. Banerjee, "An Evaluation of Multiple-Disk I/O Systems," *IEEE Transactions on Computers*, Vol. 38 (12), 1989.

[Reddy91] A. Reddy and P. Bannerjee, "Gracefully Degradable Disk Arrays," *Proceedings of FTCS-21*, 1991.

[Rudeseal92] A. Rudeseal, private communication, 1992.

[Salem86] K. Salem, H. Garcia-Molina, "Disk Striping", *Proceedings of the 2nd IEEE Conference on Data Engineering*, 1986.

[Schulze89] M. Schulze, G. Gibson, R. Katz, and D. Patterson, "How Reliable is a RAID?", *Proceedings of COMPCON*, 1989.

[TPCA89] *The TPC-A Benchmark: A Standard Specification*, Transaction Processing Performance Council, 1989.

## Appendix: Block Designs

We used six block designs in our simulations, corresponding approximately to  $\alpha$  ranging from 0.1 to 0.85. There was one case where we could not find an incomplete block design to match our requirements, and so we used a complete design. The designs below are given in the abbreviated notation described by Hall [Hall86]. In each case we give a brief description of how to derive the design from the notation, and refer the reader to Hall for more details.

In the first four designs, a set of blocks and a modulo number  $N$  are specified. The full design is generated by adding element-wise all of the residues modulo  $N$  to each of the indicated blocks. Where a period  $P$  is specified, the addition terminates after  $P$  iterations. In all cases, the addition is done modulo  $N$ .

**Block Design 1:**  $b = 70, v = 21, k = 3, r = 10, \lambda = 1, \Rightarrow \alpha = 0.1$

[0, 1, 3]; [0, 4, 10]; [0, 16, 19] (mod 21)

[0, 7, 14] (mod 21) period 7

**Block Design 2:**  $b = 105, v = 21, k = 4, r = 20, \lambda = 3 \Rightarrow \alpha = 0.15$

[0, 2, 3, 7]; [0, 3, 5, 9]; [0, 1, 7, 11]; [0, 2, 8, 11]; [0, 1, 9, 14] (mod 21)

**Block Design 3:**  $b = 21, v = 21, k = 5, r = 5, \lambda = 1 \Rightarrow \alpha = 0.2$

[3, 6, 7, 12, 14] (mod 21)

**Block Design 4:**  $b = 42, v = 21, k = 6, r = 12, \lambda = 3 \Rightarrow \alpha = 0.25$

$$[0, 2, 10, 15, 19, 20]; [0, 3, 7, 9, 10, 16] \pmod{21}$$

A block design is said to be *symmetric* if  $b = v$  and  $k = r$ . Given a symmetric block design with a particular  $b, k$ , and  $\lambda$ , a new design with  $b' = b - 1, v' = k, k' = \lambda, r' = r - 1$ , and  $\lambda' = \lambda - 1$  can be generated by selecting one of the blocks ( $B_0$ ) and constructing new blocks  $B_0', B_1', \dots, B_{b-1}'$  such that  $B_i'$  contains the  $\lambda$  objects common to  $B_i$  and  $B_0$ . Such a design is called a *derived* design.

**Block Design 5:**  $b = 42, v = 21, k = 10, r = 20, \lambda = 9 \Rightarrow \alpha = 0.45$

Derived design of  $[0, 3, 5, 8, 9, 10, 12, 13, 14, 15, 16, 20, 22, 23, 24, 30, 34, 35, 37, 39, 40] \pmod{43}$

**Block Design 6:**  $b = 1330, v = 21, k = 18, r = 1140, \lambda = 969 \Rightarrow \alpha = 0.85$

We used a complete block design in this case.